

Contributor: *Luka Anicin*

Data Scientist: *Ruslan S.*

CAPSTONE FINAL REPORT: Sentiment Analyses

PROBLEM IDENTIFICATION: Use deep learning techniques to predict sentiment on the Sentiment140 dataset.

APPROACH: Using process stages as wrangling, EDA, pre-processing, modeling, evaluating, tuning, saving.

A significant part of this project will be spent on careful data pre-processing. It includes but is not limited to retrieving related to the task data from trusted sources (like Kaggle or Tweeter websites). Then, clean it (using the wrangling process) so we can apply 'special' libraries (tools) to parse and remove noise accordingly. Next, we will gather insights that could be well applied to our models. Because we know, if data is well processed, usually it boosts the accuracy of prediction by a significant amount. Next, we split data into training and testing parts so we can use the training part to train our model and testing one to evaluate it. We will be using CNN architecture models of neural networks for dealing with sequential data (1 dimensional). Once the models are training and evaluated, we will choose the most promising one based on our requirements for further production stages (deployment part).

DATA STRUCTURE: Data is taken from Kaggle.com. It contains 1,600,000 tweets extracted using the twitter API. The tweets have been annotated (0 = negative, 1 = positive) and they can be used to detect sentiment.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1600000 entries, 0 to 1599999
Data columns (total 6 columns):
#   Column  Non-Null Count  Dtype
---  -
0   0        1600000 non-null  int64
1   1        1600000 non-null  int64
2   2        1600000 non-null  object
3   3        1600000 non-null  object
4   4        1600000 non-null  object
5   5        1600000 non-null  object
dtypes: int64(2), object(4)
memory usage: 73.2+ MB
```

Size of the data file is 227 MB (238,803,811 bytes)

- Shape is (1600000, 6)

DATA TIMEFRAME: 06/04/2009-29/05/2009

The data we were using for this project was pretty much outdated. So, if we are planning to build the model for today's production, then it will be a good practice to find more relevant data (new one).

MODEL REPONSE: Binary 0 or 1. Negative = 0, Positive = 1.

The data we were using had a pre-defined class. So, we were building the model that produces the answer in terms of two classes.

TYPE OF MODEL: Supervised Classification using Deep Learning architecture.

This project was focused on building the model using the Deep Learning approach. The most significant factor for choosing CNN was the amount of data we are going to use to train our model.

DELIVERABLES WILL BE GENERATED: PDF key points from data exploration to best model results.

DATA PRE-PROCESSING STEPS: Wrangling.

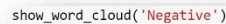
- Loading (examining the type of columns and renaming them accordingly):
 - o Column names: 'target', 'id', 'date', 'query', 'user_name', 'predictor'

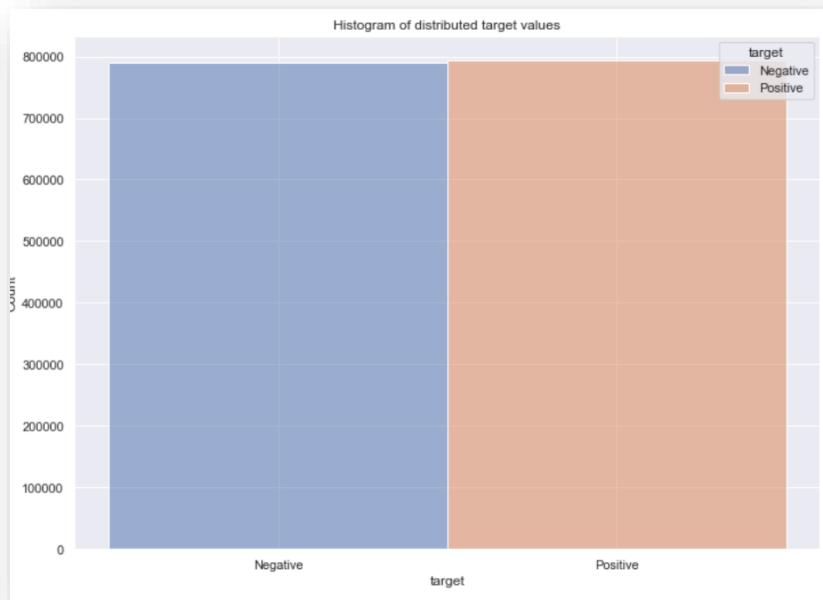
	target	id	date	query	\
0	0	1467810369	Mon Apr 06 22:19:45 PDT 2009	NO_QUERY	
1	0	1467810672	Mon Apr 06 22:19:49 PDT 2009	NO_QUERY	
2	0	1467810917	Mon Apr 06 22:19:53 PDT 2009	NO_QUERY	
	user_name				predictor
0	_TheSpecialOne_	@switchfoot	http://twitpic.com/2y1zl	- Awww, t...	
1	scotthamilton		is upset that he can't update his Facebook by ...		
2	mattycus	@Kenichan	I dived many times for the ball. Man...		

- Dropping (identifying dependent and independent variables, irrelevant are dropped):
 - o Dropped columns: 'id', 'date', 'query', 'user_name'

- EDA:** Visualization part to get insights.

- ```
show_word_cloud('Positive')
```





**FEATURES:** There is only one column that will be used as independent feature to train model.

- Feature column name: 'predictor'

|         | target   | predictor                                         |
|---------|----------|---------------------------------------------------|
| 587328  | Negative | tried still work suggestions                      |
| 433711  | Negative | photo idea also back amp really itchy             |
| 1058707 | Positive | filw maritza 6121 hey flwrs wl try thank ya m...  |
| 211457  | Negative | pity djokovic                                     |
| 1414244 | Positive | head 7 hours driving group people numerous win... |

**MODEL DESCRIPTION:** Architecture and compilation parts.

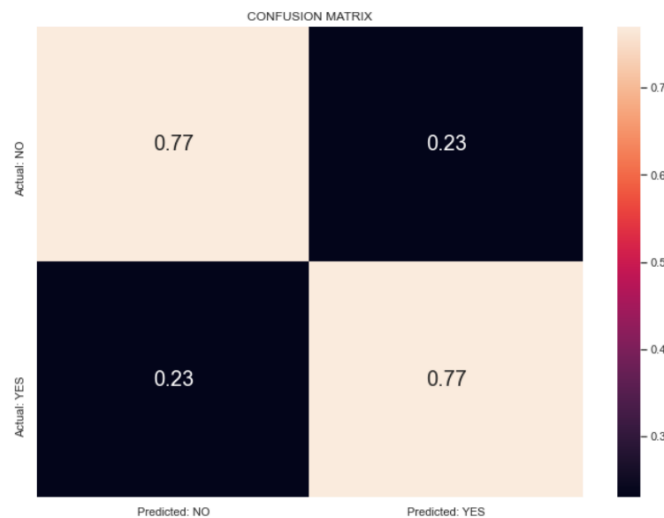
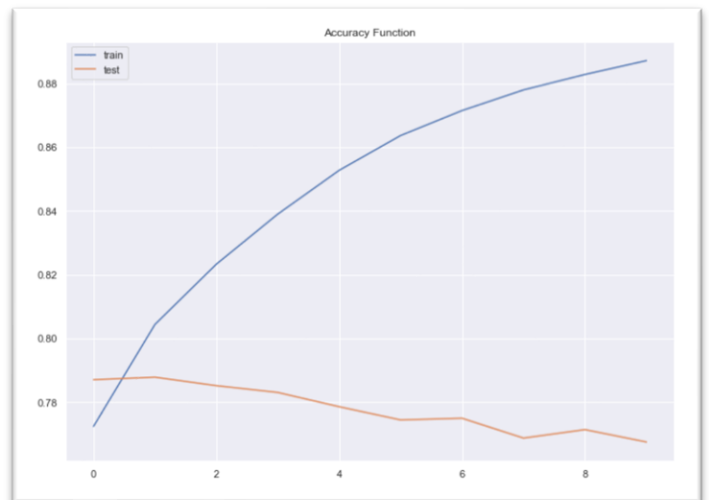
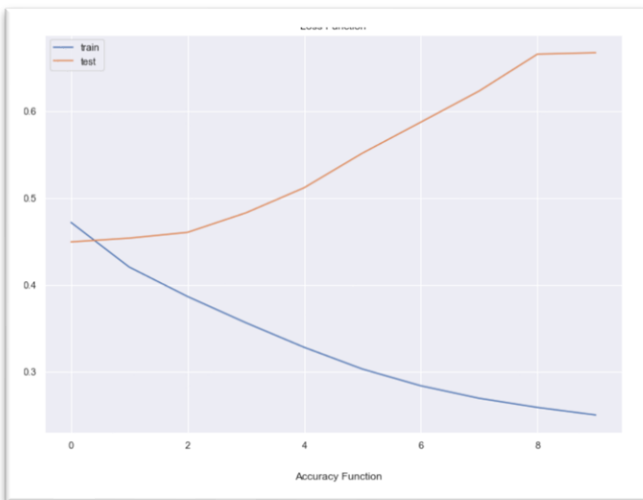
- Model is 'Sequential' (CNN convolutional 1 dimension) that consists of 5 hidden layers
- Embedding layer:
  - o Input data size (length) is 290412
  - o Maximum sequence length is 30
  - o Embedding dimension is 100
- SpatialDropout1D layer:
  - o Fraction of the input units to drop is 0.2
- LSTM (long-sort term memory) layer:
  - o Units (positive integer, dimensionality of the output space) is 100

- Dropout (fraction of the units to drop for the linear transformation of the inputs) is 0.2
- Recurrent dropout (fraction of the units to drop for the linear transformation of the recurrent state) is 0.2
- Dense layer:
  - Units (positive integer, dimensionality of the output space) is 15
  - Activation (activation function to use) is 'relu' (the rectified linear activation function)
- Dense layer:
  - Units (positive integer, dimensionality of the output space) is 1
  - Activation (activation function to use) is 'sigmoid' (the rectified linear activation function)
- Model's compiler:
  - Optimizer (name of optimizer) is 'adam'
  - Loss (name of loss function) is 'binary\_crossentropy'
  - Metrics (list of metrics to be evaluated by the model during training and testing) is 'accuracy'
- Epoch size is 10
- Batch size is 512
- Total parameters: 29,123,131
- Trainable parameters: 29,123,131
- Non-trainable parameters: 0

| Layer (type)                 | Output Shape    | Param #  |
|------------------------------|-----------------|----------|
| embedding_3 (Embedding)      | (None, 30, 100) | 29041200 |
| spatial_dropout1d_2 (Spatial | (None, 30, 100) | 0        |
| lstm_2 (LSTM)                | (None, 100)     | 80400    |
| dense_7 (Dense)              | (None, 15)      | 1515     |
| dense_8 (Dense)              | (None, 1)       | 16       |
| Total params: 29,123,131     |                 |          |
| Trainable params: 29,123,131 |                 |          |
| Non-trainable params: 0      |                 |          |

**MODEL PERFORMANCE:** Chosen (final) model is number 4.

|                     | precision | recall | f1-score | support |
|---------------------|-----------|--------|----------|---------|
| Class 0 - Negatives | 0.77      | 0.77   | 0.77     | 157487  |
| Class 1 - Positives | 0.77      | 0.77   | 0.77     | 159252  |
| accuracy            |           |        | 0.77     | 316739  |
| macro avg           | 0.77      | 0.77   | 0.77     | 316739  |
| weighted avg        | 0.77      | 0.77   | 0.77     | 316739  |



**PERFORMENS OF ALL TRAINED MODELS:** Metrics of all trained models.

| Model_name | Precision | Recall | F1-score | Accuracy |
|------------|-----------|--------|----------|----------|
| Model_1    | 0.75      | 0.75   | 0.75     | 0.75     |
| Model_3    | 0.75      | 0.75   | 0.75     | 0.75     |
| Model_4    | 0.77      | 0.77   | 0.77     | 0.77     |

**MODEL FINDINGS:** Predictor feature (tweets column) is the only column that was used to train our model. It is significant to our task.

Because our data were well balanced our model made good predictions toward true positives and true negatives.

**NEXT STEPS:** The model is ready for production. It also is ready for the further hyper-parameter improvements and new data fitting.

**IDEAS FOR FURTHER RESEARCH:** There are some steps that could be done in case to improve overall performance of this model.

First of all, we could continue tuning our final model by using different hyper-parameters. Unfortunately, it takes a significant amount of time and proper hardware resources, which I do not possess at this moment. Also, we can apply additional libraries dealing with text. For example, we can fix words that are grammatically not correct (so, we can use word embedding with pre-trained weights more effectively). Also, we can define and split words that were used as a combination of two or more for further sentiment analyses. Finally, we could play with the architecture of our model by creating more hidden layers in case to reach the desired point. But we should stay cautious by not making the model too complex and overfitted.

**CREDITS:** A huge thanks to Luka Anicin for being an amazing teacher and Springboard mentor.

## **CITATIONS:**

Team, K. (n.d.). *Keras documentation: LSTM layer*. Keras. Retrieved September 21, 2021, from [https://keras.io/api/layers/recurrent\\_layers/lstm/](https://keras.io/api/layers/recurrent_layers/lstm/).

Team, K. (n.d.). *Keras documentation: Dense layer*. Keras. Retrieved September 21, 2021, from [https://keras.io/api/layers/core\\_layers/dense/](https://keras.io/api/layers/core_layers/dense/).

*Tf.keras.model : tensorflow core v2.6.0*. TensorFlow. (n.d.). Retrieved September 21, 2021, from [https://www.tensorflow.org/api\\_docs/python/tf/keras/Model#compile](https://www.tensorflow.org/api_docs/python/tf/keras/Model#compile).

*Removing stop words with NLTK in python*. GeeksforGeeks. (2021, May 31). Retrieved September 21, 2021, from <https://www.geeksforgeeks.org/removing-stop-words-nltk-python/>.

Go, A., Bhayani, R. and Huang, L., 2009. Twitter sentiment classification using distant supervision. CS224N Project Report, Stanford, 1(2009), p.12.